

# Hyperparameter Learning for Deep Learning-based Recommender Systems

Di Wu, *Member, IEEE*, Bo Sun, and Mingsheng Shang

**Abstract**—Deep learning (DL)-based recommender system (RS), particularly for its advances in the recent five years, has been startling. It reshapes the architectures of traditional RSs by lifting their limitations in dealing with data sparsity and cold-start issues. Yet, the performance of DL-based RS, like many other DL-based intelligent systems, heavily relies on selecting hyperparameters. Unfortunately, the most common selection approach is still Grid Search that requires numerous computational resources and human efforts. Motivated by this, this paper proposes a general hyperparameter optimization framework, named DE-Opt, which can be applied to most existing DL-based RSs seamlessly. The main idea of DE-Opt is to incorporate differential evolution (DE) into the model training process of a DL-based RS at a layer-wise granularity to simultaneously auto-learn its two key hyperparameters, namely, the learning rate  $\eta$  and the regularization coefficient  $\lambda$ . As a result, the system performance in terms of recommendation accuracy and computational efficiency is uplifted. Experimental results on three benchmark datasets substantiate that: 1) DE-Opt is compatible with the most recent advancements of DL-based RS to automate their hyperparameter-tuning processes, and 2) DE-Opt excels among its state-of-the-art hyperparameter-optimization competitors in terms of both higher learning performance and lower runtime.

**Index Terms**—Online services, recommender systems, hyperparameter tuning, grid search, deep learning, differential evolution.

## 1 INTRODUCTION

RECOMMENDER SYSTEM (RS) has become an indispensable building block in online services for business-boosting and user-experience elevation [1, 2]. Traditional RS tends to yield inferior performance in resolving the data sparsity and cold-start problems without sacrificing recommendation quality [1, 2, 4]. Deep learning (DL) has been advocated as a plausible solution to the issue, thanks to its strong capability of modeling complex and non-linear user-item interactions [5, 6]. In recent DL-based RSs, superior recommendation performance is observed [3, 5, 33-34].

Despite the superiority, existing DL-based RSs mainly suffer from a common drawback—the tedious process of hyperparameter tuning [7, 34]. On the one hand, due to the black-box nature of DL model architectures, the hyperparameters in a DL-based RS directly decide the system properties in both training and trained states, i.e., how fast the RS can converge and how well the converged equilibrium optimizes the designed learning objective [38-39]. On the other hand, the ubiquity of DL-based RSs exacerbates this nature. Various regularization terms have been engineered to incorporate a diversity of side information from the items and users (e.g., textual description of items and geographical relations among users) in a wide range of recommendation applications. Such side information is the key to making the DL-based RS great [7-9, 34, 38-39]. While, each regularization term correspondingly added to the learning objective introduces a new hyperparameter, making the tuning effort even more cumbersome. A concrete example that substantiates the impact of hyperparameter tuning in a DL-based RS is presented as follows.

- D. Wu is with the College of Computer and Information Science, Southwest University, Chongqing 400715, China, and also with the Chongqing Institute of Green and Intelligent Technology, Chinese Academy of Sciences, Chongqing 400714, China (email: wudi.cigit@gmail.com).
- B. Sun and M. Shang are with the Chongqing Key Laboratory of Big Data and Intelligent Computing, Chongqing Institute of Green and Intelligent Technology, Chinese Academy of Sciences, Chongqing 400714, China; and also with the Chongqing School, University of Chinese Academy of Sciences, Chongqing 400714, China (e-mail: {sunbo, ms Shang}@cigit.ac.cn).
- Corresponding Author: Mingsheng Shang.

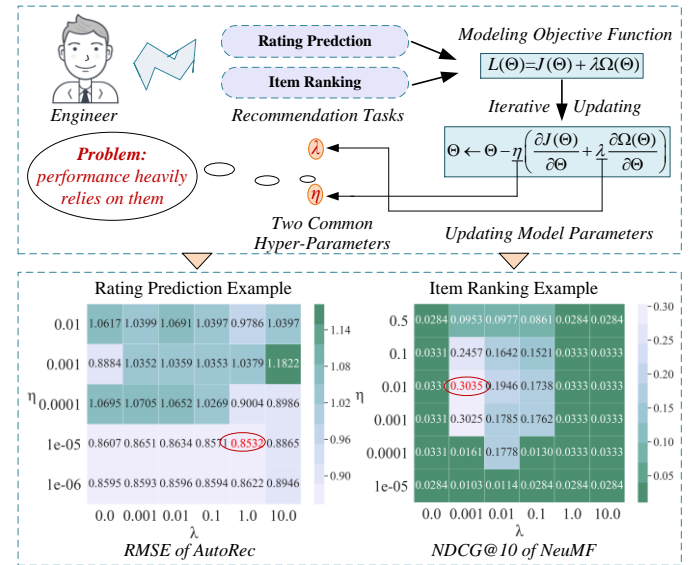


Fig. 1. An example of hyperparameters tuning: AutoRec and NeuMF are highly sensitive to  $\eta$  and  $\lambda$  (lower RMSE and higher NDCG@10 are better).

ons among users) in a wide range of recommendation applications. Such side information is the key to making the DL-based RS great [7-9, 34, 38-39]. While, each regularization term correspondingly added to the learning objective introduces a new hyperparameter, making the tuning effort even more cumbersome. A concrete example that substantiates the impact of hyperparameter tuning in a DL-based RS is presented as follows.

**Example 1.** A typical DL-based recommendation model is trained by optimizing an objective taking the form of  $L(\Theta) = J(\Theta) + \lambda\Omega(\Theta)$  [5], where  $\Theta$  parameterizes the RS and  $J(\Theta)$  and  $\Omega(\Theta)$  denote the loss and regularization terms, respectively. Two hyperparameters exist in this example, as shown in the top panel of Figure 1. The first is  $\lambda$  in the learning objective that

balances the scale between the loss and regularization terms. The second is the learning rate, denoted by  $\eta$ , which decides how fast the RS is optimized based on stochastic oracles (e.g., gradients or Hessians) in an iterative fashion [15, 16]. To test the significance of  $\lambda$  and  $\eta$ , we hereby evaluate two representative DL-based RS models, i.e., AutoRec [22] (rating prediction) and NeuMF [14] (item ranking), on the MovieLens 1M dataset downloaded from MovieLens<sup>1</sup>. Their performance is presented in the bottom panel of Figure 1. We observe that the variation of  $\lambda$  and  $\eta$  leads to the extensive changes of model performance, where AutoRec and NeuMF respectively achieve different root mean squared error (RMSE) [22] and normalized discounted cumulative gain (NDCG, NDCG@10 denotes top 10) [14]. The lighter the color, the better the performance. We note that the optimal models are resulted from different values of  $\lambda$  and  $\eta$  (colored in red).

Example 1 reveals how the performance of a DL-based RS is heavily dependent on an ad-hoc setting of hyperparameters, which has been manifested in a variety of DL-based RS designs [7-9, 35]. Thus far, to tune the RS hyperparameters, the most common approach is still Grid Search [10, 15], where the optimal hyperparameters are searched in predefined intervals. This search inevitably incurs a combinatorial complexity, which can consume numerous computational resources and human efforts. Therefore, alternatives to Grid Search are desired for efficient and effective hyperparameter searching.

Motivated by this, this paper proposes a general hyperparameter optimization framework for DL-based RSs based on differential evolution (DE) [18-20], named DE-Opt. Its key ideas are two-fold as follows. 1) DE-Opt decomposes and controls the RS training process at a finer level of granularity by allowing the hyperparameters to differ across neural layers. This enlarges the search space for optimal hyperparameters yet alleviates the impact of each hyperparameter so as to better the resultant models. 2) DE-Opt frames the hyperparameter tuning task as an optimization problem and embeds it into the RS training framework, so that the model parameters and the hyperparameters of an RS are jointly trained. As such, DE-Opt allows to train the RSs only once (in contrast, Grid Search requires to train RSs repeatedly by trying a different hyperparameter combination per each repeat), thereby substantially improving the efficiency and effectiveness of DL-based RSs by attaining their optimal RS models and hyperparameters all at once.

The specific contributions of this paper include:

- This is the first work to propose the auto-learning of the hyperparameters of DL-based RSs at a layer-wise granularity, where DE is employed to solve the nonconvex and combinatorial optimization problems.
- The proposed DE-Opt framework is agnostic to the learning objectives and thus is compatible with most existing DL-based RSs, boosting their training effectiveness and efficiency by making their hyperparameters auto-learned along with the model parameters jointly.

- Theoretical analyses and algorithm designs are provided for the proposed DE-Opt framework.

Extensive experiments on three real-world datasets are conducted. The results verify that our DE-Opt: 1) can significantly improve the state-of-the-art DL-based RSs by making their hyperparameters be auto-learned, and 2) significantly outperforms the state-of-the-art adaptive hyperparameter tuning competitors. Besides, to promote reproducible research, we open-access our code at the following link: <https://github.com/Wuziqiao/DE-Opt>.

## 2 RELATED WORK

In a nutshell, the mainstream RS designs can be categorized into two settings [5, 15, 16], i.e., **rating prediction** [36-37] and **item ranking** (i.e., top-K recommendation) [14]. The pioneering study that implements deep neural networks into rating prediction task is AutoRec [22], wherein an autoencoder architecture realizes collaborative filtering [46]. Since then, a flurry of deep learning (DL)-based rating prediction RS has been proposed. In [23], a rating prediction model with abstractive tips generation (NRT) is devised by combining multilayer perceptron and recurrent neural networks. More recently, [24] proposes a federated rating prediction model (MetaMF) based on meta-learning [38], and [28] proposes a multi-component graph convolutional collaborative filtering method for rating prediction. For item ranking, a classic model is NeuMF [14] which firstly combines matrix factorization and multilayer perceptron. Later, [25] proposes a latent relevant metric learning model (LRML) to learn the correlation between users and items. Graph convolution network (GCN) is also employed to implement item ranking RSs. Proposed by He et al., two state-of-the-art GCN-based RSs, namely, NGCF [21] and LightGCN [27], have been devised to achieve remarkable item ranking recommendation performance.

Nevertheless, these methods all suffer a limitation that their hyperparameters  $\lambda$  and  $\eta$  must be chosen in prior. To be applied in practice, a process of hyperparameter tuning, such as Grid Search [10, 15], has to be tortured. Typically, a combination of multiple candidate hyperparameter sets needs to be tried repeatedly for the RS model. The best set of hyperparameters yielding from a held-out validation dataset is selected. This is a tedious effort that consumes considerable time, manpower, and computational resources.

To alleviate this limitation, several efforts are made to automate the hyperparameter tuning of RSs. A classic study is SGDA [11] which updates the regularization coefficient  $\lambda$  and the RS model parameters via alternating optimization. Inspired by SGDA,  $\lambda$ Opt is proposed [8] to tune  $\lambda$  adaptively by enforcing regularization during training, which also reveals that an adaptive regularization at a layer-wise finer level can bring performance benefits. As a parallel effort, several studies focused on tuning the other significant hyperparameter -- the learning rate  $\eta$ . A representative  $\eta$ -adaptive optimizer is Adam [17], which is a first-order gradient-based stochastic optimizer

<sup>1</sup> <https://grouplens.org/datasets/movielens/>

based on adaptive estimation of lower-order momentum. Later, AMSGrad [12] and AdaMod [13] were proposed to aid its sub-optimal convergence rate and the ill-conditioned initialization of Adam by bounding the regularities of momenta.

Unfortunately, these approaches can only tune a single hyperparameter  $\lambda$  or  $\eta$  at each time, while no one works for both. As such, if one covets to simultaneously tune both  $\lambda$  and  $\eta$  adaptively at once, a compromise method is yielded. For example, a straightforward adaptation is to employ an  $\eta$ -adaptive optimizer (e.g., Adam) to optimize a  $\lambda$ -adaptive model (e.g.,  $\lambda$ Opt). Such a method, alas, boils down to be a slightly improved Grid Search, leading to the suboptimal pairing of hyperparameter set and the model parameters and ending with inferior learning efficiency and performance.

To counter the issue, this paper presents DE-Opt, which significantly differs from the prior studies in the sense that it can make existing DL-based RSs'  $\lambda$  and  $\eta$  be simultaneously auto-learned at layer-granularity (i.e., each layer is allowed to be equipped with different  $\lambda$  and  $\eta$  pairs) during the training on-the-fly. Moreover, the proposed DE-Opt is a generic framework that is compatible with most existing DL-based RSs to boost their learning performance and training efficiency.

### 3 PRELIMINARIES

#### 3.1 Symbols and Notations

For clarification, we summarize the main notations employed in Table 1. The complete symbols and notations are deferred to Table S.1 in the Supplementary File.

#### 3.2 General Learning Paradigm of DL-based RS

Despite its multiple variants, the learning procedure of the DL-based RS can be uniformly framed as follows [24]:

$$\begin{cases} l^1 = f(W^1x + b^1), \\ l^n = f(W^n l^{n-1} + b^n), \\ l^N = f(W^N l^{N-1} + b^N), \end{cases} \quad (1)$$

where  $x$  indicates the initialized (user or item) embedding vector,  $l^n$  indicates the output vector of  $n$ -th layer after activation,  $f(\cdot)$  indicates the activation function,  $W^n$  indicates the  $n$ -th weight matrix,  $b^n$  denotes the  $n$ -th bias vector,  $n \in \{2, \dots, N\}$ , and  $l^N$  indicates the final output of the last layer after activation.

To model a DL-based RS, the design of the objective function  $L(\Theta) = J(\Theta) + \lambda\Omega(\Theta)$  is crucial [5, 40-41]. Specifically, point-wise loss and pair-wise loss are commonly used to formulate the loss term  $J(\Theta)$  of rating prediction and item ranking, respectively [5, 15]. For **rating prediction**, a general objective function can be usually formulated with square distance as follows:

$$L_{\text{Rating}}(\Theta) = \sum_{(u,j) \in Z} \left( (\hat{y}_{u,j} - y_{u,j})^2 | \Theta \right) + \lambda\Omega(\Theta), \quad (2)$$

where  $y_{u,j}$  denotes the observed interaction between a user  $u$  and an item  $j$ . Denoted by  $\hat{y}_{u,j}$  is the prediction of the ground-truth  $y_{u,j}$ , and by  $Z$  is the observed explicit rating set of the input matrix which is sparse and high-dimensional in its argument.

For **item ranking**, (2) can not be used for binarized implicit

TABLE 1. SYMBOLS AND NOTATIONS.

Symbol	Explanation
$L(\cdot), \Theta$	The objective function and its model parameters, respectively.
$J(\cdot), \Omega(\cdot)$	The loss and regularization terms of $L(\cdot)$ , respectively.
$L_{u,j}(\cdot), J_{u,j}(\cdot)$	The state of $L(\cdot)$ and $J(\cdot)$ on $y_{u,j}$ , respectively.
$W, b$	The weight matrix and bias vector sets of $N$ layers, respectively.
$l^n$	The output vector of $n$ -th layer after activation, $n \in \{1, 2, \dots, N\}$ .
$W^n, b^n$	The weight matrix and bias vector of $n$ -th layer, respectively.
$x$	Initializing embedding vector.
$f(\cdot)$	The activation function.
$Z, D$	The observed explicit ratings and training triplets sets of input sparse matrix for rating prediction and item ranking, respectively.
$y_{u,j}, \hat{y}_{u,j}$	The observed ground truth of interaction between user $u$ and item $j$ and its corresponding prediction.
$\lambda^n, \eta^n$	The regularization coefficient and learning rate of $n$ -th layer, respectively.
$\Delta^n$	The term <i>w.r.t.</i> partial derivative at $n$ -th layer.
$H, \Lambda$	The fine-grained vectorized $\eta$ and $\lambda$ of $N$ layers, respectively.
$NP$	The total number of individuals of DE.
$X_i$	The $i$ -th target vector (individual), $i \in \{1, 2, \dots, NP\}$ .
$V_i$	The mutant vector generated for each target vector $X_i$ .
$r_1, r_2, r_3$	The randomly chosen number from $\{1, 2, \dots, NP\}$ .
$F_i$	The scaling factor controlling the scaling of various vectors.
$U_i$	The new trial vector generated from $X_i$ and $V_i$ .
$\text{perf}(\cdot)$	The fitness/performance function.
$\Gamma, \Psi$	The testing set and its contained users set, respectively.
$K$	Cutoff numbers of items in the given ranked list to each user.
$\Gamma_K(u)$	The intersection set between the top- $K$ items generated by the tested model and the testing items of $\Gamma$ by user $u$ .

feedback well [26]. Alternatively, the pair-wise with Bayesian Personalized Ranking (BPR) loss [26] is commonly adopted. Its general objective function is formulated as follows:

$$L_{\text{Ranking}}(\Theta) = - \sum_{(u,j,j') \in D} \ln \left( f \left( (\hat{y}_{u,j} - \hat{y}_{u,j'}) | \Theta \right) \right) + \lambda\Omega(\Theta), \quad (3)$$

where  $D = \{(u, j, j') | y_{u,j}=1, y_{u,j'}=0\}$  denotes the training triplets set,  $\hat{y}_{u,j}$  and  $\hat{y}_{u,j'}$  denote the predictions of interactions  $(u, j)$  and  $(u, j')$ , respectively.

Finally, the training of DL-based recommendation models is to minimize their objective function [27]. To this end, stochastic gradient descent (SGD) [28-30] is most commonly adopted.

## 4 THE PROPOSED DE-OPT OPTIMIZATION FRAMEWORK

### 4.1 Motivation of DE-Opt

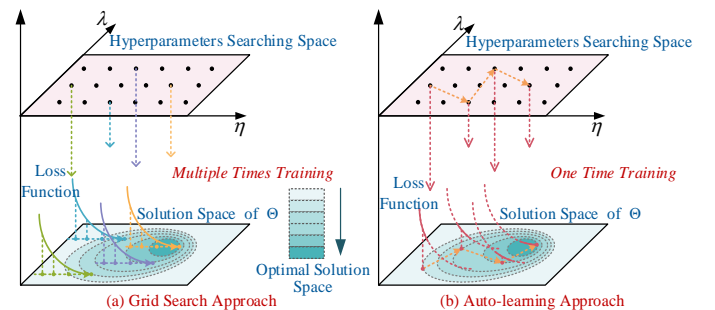


Fig. 2. The differences of hyperparameters tuning between Grid Search and auto-learning approach.



The performance of a DL-based recommendation model heavily relies on hyperparameters tuning  $\eta$  and  $\lambda$ . Figure 2 visualizes the differences between the two tuning approaches of Grid Search and auto-learning, where each one has a hyperparameter searching space (candidate values of  $\eta$  and  $\lambda$ ) and solution space of  $\Theta$  (desired model parameters  $\Theta$ ), and the arrows represent the searching paths.

**Grid Search.** First,  $\eta$  and  $\lambda$  are initialized with two values from the searching space. After that, they become a pair of fixed values in all the steps of minimizing the loss function. Each pair of  $\eta$  and  $\lambda$  in the searching space corresponds to a solution in the solution space. Figure 2 (a) shows that if we want to find the optimal solution of  $\Theta$ , we need to traverse the complete searching space where each search requires full training steps. As a result, Grid Search either consumes computational resources or requires prior knowledge about selecting appropriate  $\eta$  and  $\lambda$  candidates.

**Auto-learning.** Different from Grid Search, auto-learning allows employing different  $\eta$  and  $\lambda$  at each training step. Figure 2 (b) shows that auto-learning only needs one-time training to search for the optimal  $\Theta$ . Although auto-learned trajectories can greatly reduce computational consumption, there are two key observations regarding the deficiencies of this auto-learning process, which motivate our DE-Opt design as follows.

- At different training stages, the learning rate of gradient descent and strength of regularization should be different. Since the model has not learned much from data at the early training stages, the learning rate should be large and the regularization should be light. In contrast, a small learning rate and strong regularization are required to finely search optimal  $\Theta$  and avoid overfitting at the end training stages.
- Each layer should have a different learning rate of gradient descent and strength of regularization. During the backpropagation of errors, the gradient may vanish or explode from the last layer to the first layer, preventing to search for the globally optimal  $\Theta$ . One possible way of alleviating this problem is to set different  $\eta$  and  $\lambda$  at layer-granularity.

## 4.2 The Architecture of DE-Opt

Upon the two observations, we aim to build an auto-learning scheme for DL-based recommendation models'  $\eta$  and  $\lambda$  at layer granularity simultaneously during the training process. Note that this auto-learning problem is discrete. To solve it, we employ differential evolution (DE) [18-20], which is known for its high efficiency and robustness in obtaining global optima in discrete numerical optimization tasks. In particular, our DE-Opt architectural design, as shown in Figure 3, consists of two parts. Part I shows how to train a DL-based recommendation model layer by layer, where  $\lambda$  and  $\eta$  are designed at layer-granularity, i.e., each layer has different  $\lambda$  and  $\eta$ . Part II is to employ DE to simultaneously auto-learn each layer's  $\lambda$  and  $\eta$  at each epoch, including four operations, i.e., initialization, mutation, crossover, and selection. We present the details of the two building blocks in the following sections.

### 4.2.1 Part I: Training a DL-based Recommendation Model with Layer-wise Granular Hyperparameter Control

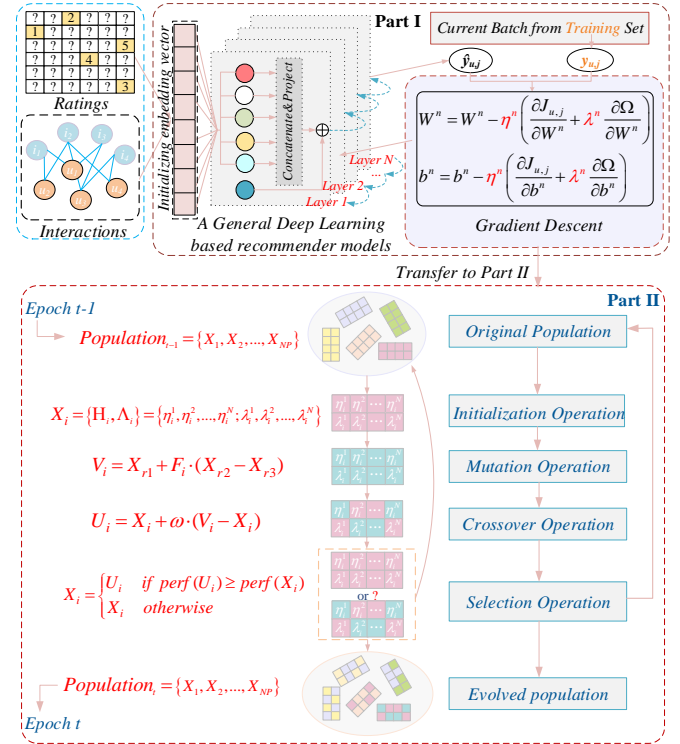


Fig. 3. The Architecture of DE-Opt.

TABLE 2. THE VECTORIZED HYPERPARAMETERS OF  $\lambda$  AND  $\eta$ .

Hyperparameters	Layer 1	...	Layer $n$	...	Layer $N$
Regularization Coefficient: $\Lambda$	$\lambda^1$		$\lambda^n$		$\lambda^N$
Learning Rate: $H$	$\eta^1$		$\eta^n$		$\eta^N$

To finely control different layers,  $\lambda$  and  $\eta$  are matched with the number of hidden layers, i.e., model depth, as shown in Table 2. The  $\Lambda$  and  $H$  denote the  $\lambda$  and  $\eta$  vectors, respectively, where  $\lambda^n$  and  $\eta^n$  denote  $\lambda$  and  $\eta$  at  $n$ -th layer,  $n \in \{1, 2, \dots, N\}$ , respectively.

Supposing input the observed ground-truth  $y_{u,j}$  of interaction  $(u, j)$  between user  $u$  and item  $j$  at the current batch in the training set, we have the current objective function as follows:

$$L_{u,j}(\Theta) = L_{u,j}(W, b) = J_{u,j}(W, b) + \Lambda \square \Omega(W, b), \quad (4)$$

where  $J_{u,j}(\cdot)$  indicates the loss function on  $y_{u,j}$ ,  $W$  is the weight matrix sets of  $N$  layers,  $b$  is the bias vector sets of  $N$  layers, and  $\square$  indicates the element-wise product between parameters of each layer and  $\Lambda$ . By employing SGD to train (4) [15-17], we obtain the training rules for  $n$ -th layer's model parameters as follows:

$$W^n \leftarrow W^n - \eta^n \frac{\partial L_{u,j}}{\partial W^n} = W^n - \eta^n \left( \frac{\partial J_{u,j}}{\partial W^n} + \lambda^n \frac{\partial \Omega}{\partial W^n} \right), \quad (5)$$

$$b^n \leftarrow b^n - \eta^n \frac{\partial L_{u,j}}{\partial b^n} = b^n - \eta^n \left( \frac{\partial J_{u,j}}{\partial b^n} + \lambda^n \frac{\partial \Omega}{\partial b^n} \right). \quad (6)$$

Note that in (5) and (6), the terms of  $\partial \Omega / \partial W^n$  and  $\partial \Omega / \partial b^n$  can be directly computed while that of  $\partial J_{u,j} / \partial W^n$ , and  $\partial J_{u,j} / \partial b^n$  require the chain rule [5-7] to compute. Then, we can compute  $\partial J_{u,j} / \partial W^n$

and  $\partial J_{u,j}/\partial b^n$  as follows:

$$\begin{aligned} \frac{\partial J_{u,j}}{\partial W^n} &= \frac{\partial J_{u,j}}{\partial l^N} \frac{\partial l^N}{\partial W^n} \\ &= \frac{\partial J_{u,j}}{\partial l_N} \frac{\partial f(W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})} \frac{\partial (W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})} \dots \\ &\quad \frac{\partial f(W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \frac{\partial (W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \frac{\partial f(W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} l^{n-1} \end{aligned} \quad (7)$$

$$\begin{aligned} \frac{\partial J_{u,j}}{\partial b^n} &= \frac{\partial J_{u,j}}{\partial l^N} \frac{\partial l^N}{\partial b^n} \\ &= \frac{\partial J_{u,j}}{\partial l_N} \frac{\partial f(W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})} \frac{\partial (W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})} \dots \\ &\quad \frac{\partial f(W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \frac{\partial (W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \frac{\partial f(W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \end{aligned} \quad (8)$$

Let  $\Delta^n$  be the term *w.r.t.* partial derivative at  $n$ -th layer as follows:

$$\Delta^n = \frac{\partial (W^{n+1} l^n + b^{n+1})}{\partial f(W^{n+1} l^n + b^{n+1})} \frac{\partial f(W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} \frac{\partial (W^{n+1} l^n + b^{n+1})}{\partial (W^{n+1} l^n + b^{n+1})} l_{n-1}^{n-1}. \quad (9)$$

Then,  $\partial J_{u,j}/\partial W^n$  and  $\partial J_{u,j}/\partial b^n$  can be computed by error backpropagation from the last  $N$ -th layer to the  $n$ -th layer as follows [5-7]:

$$\begin{aligned} \frac{\partial J_{u,j}}{\partial W^n} &= \frac{\partial J_{u,j}}{\partial l^N} \frac{\partial l^N}{\partial W^n}, & \frac{\partial J_{u,j}}{\partial b^n} &= \frac{\partial J_{u,j}}{\partial l^N} \frac{\partial l^N}{\partial b^n}, \\ \frac{\partial l^N}{\partial W^n} &= \frac{\partial f(W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})} l^{N-1}, & \frac{\partial l^N}{\partial b^n} &= \frac{\partial f(W^N l^{N-1} + b^N)}{\partial (W^{N-1} l^{N-2} + b^{N-1})}, \\ \frac{\partial l^N}{\partial W^n} &= \frac{\partial l^N}{\partial W^{n+1}} \Delta^n, & \frac{\partial l^N}{\partial b^n} &= \frac{\partial l^N}{\partial b^{n+1}} \Delta^n \frac{l_n}{l_{n-1}}, \end{aligned} \quad (10)$$

Therefore, we can finely control training process of each layer with different  $\eta$  and  $\lambda$  at layer-granularity based on (5)–(10).

#### 4.2.2 Part II: DE-Based Hyperparameter Auto-Learning

This part adopts DE to auto-learn  $\Lambda$  and  $H$  simultaneously, including Initialization, Mutation, Crossover, and Selection [18-20, 42-45].

**Initialization.** At  $(t-1)$ -th training epoch, it starts from initializing a pair of  $\Lambda$  and  $H$  with  $NP$  individuals, where each pair is denoted as a ‘target vector’ as follows:

$$X_i = \{H_i, \Lambda_i\} = \left\{ \left[ \eta_i^1, \eta_i^2, \dots, \eta_i^N \right], \left[ \lambda_i^1, \lambda_i^2, \dots, \lambda_i^N \right] \right\}, \quad (11)$$

where  $X_i$  denotes the  $i$ -th individual,  $i \in \{1, 2, \dots, NP\}$ . All the  $NP$  individuals form an original population.

**Mutation.** Each target vector  $X_i$  has a corresponding mutant vector  $V_i$ . The *DE/Rand/1* mutation strategy is adopted to generate  $V_i$  [18-20] as follows:

$$V_i = X_{r1} + F_i \cdot (X_{r2} - X_{r3}), \quad (12)$$

where  $r1, r2, r3 \in \{1, 2, \dots, NP\}$  are random and different from each other and  $i$ , and  $F_i$  is a scaling factor controlling the scaling of various vectors. Note that there are some other frequently adopted mutation strategies, including *DE/Best/1*, *DE/RandToBest/1*, *DE/Rand/2*, *DE/Best/2*, and *DE/RandToBest/2*. They are introduced in the Supplementary File. In particular,  $F_i$  controls the balance of optimal performance and convergence

speed. Previous scale factor local search in DE (SFLSDE) method [19, 20] is effective in making  $F_i$  self-adaptive as follows:

$$F_i = \begin{cases} SFGSS & \text{if } rand_3 < \tau_2 \\ SFHC & \text{if } \tau_2 \leq rand_3 < \tau_3 \\ F_i + F_u \cdot rand_1 & \text{if } rand_2 < \tau_1 \text{ and } rand_3 > \tau_3 \\ F_i + F_u \cdot rand_1 & \text{if } rand_2 \geq \tau_1 \text{ and } rand_3 > \tau_3 \end{cases}, \quad (13)$$

where  $rand_1$ ,  $rand_2$ , and  $rand_3$  are uniformly random values between 0 and 1;  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$  are constant thresholds. According to [19, 20], all the parameters in (13) are recommended as follows:  $SFGSS=8$ ,  $SFHC=20$ ,  $F_i=0.1$ ,  $F_u=0.9$ ,  $\tau_1=0.1$ ,  $\tau_2=0.03$ ,  $\tau_3=0.07$ , and  $F_i$  is initialized to a random value between 0 and 1.

**Crossover.** To increase population diversity,  $V_i$  is used to disturb  $X_i$ . Crossover operation is to generate a new trial vector  $U_i$  based on  $X_i$  and  $V_i$ . The *arithmetic crossover* strategy [18-20] is adopted to generate the trial vector  $U_i$  as follows:

$$U_i = X_i + \omega \cdot (V_i - X_i), \quad (14)$$

where  $\omega$  is a random number in  $[0, 1]$ . By combining (12) into (14), we have:

$$U_i = X_i + \omega \cdot (X_{r1} - X_i) + F_i \cdot (X_{r2} - X_{r3}). \quad (15)$$

Besides, the *binomial crossover* strategy is also usually used. It is introduced in the Supplementary File.

**Selection.** Selection operation is to test whether  $U_i$  or  $X_i$  to make a DL-based recommendation model achieve better performance on the validation dataset as follows:

$$X_i = \begin{cases} U_i & \text{if } perf(U_i) \geq perf(X_i) \\ X_i & \text{otherwise} \end{cases}, \quad (16)$$

where  $perf(\cdot)$  is a function that adopts a specific metric to evaluate a DL-based recommendation model’s performance.

Next, DE-Opt moves to  $t$ -th training epoch after all the  $NP$  individuals (target vectors) are evolved, in which  $H$  and  $\Lambda$  can be set by any one of the evolved  $NP$  target vectors. Note that with such operations of part II, the auto-learned  $H$  and  $\Lambda$  can guarantee that the performance of the next training epoch is always better than, or at least the same to, the current epoch. Therefore, DE-Opt guarantees to boost the performance of a DL-based RS.

#### 4.3 Theoretical Analysis and Proof

**Theorem 1.** Given the searching space of  $\eta$  and  $\lambda$ , DE-Opt is convergent in searching the values of  $\eta$  and  $\lambda$  to make a DL-based recommendation model perform better.

**Proofs of theorem 1.** To begin with, we introduce Markov chain of DE-Opt and convergence criteria. Then, we prove the convergence of DE-Opt.

##### 4.3.1 Markov Chain of DE-Opt

Following previous studies [29-32], we build a discrete Markov chain for DE-Opt.

**Initialization.** From formula (11), the state of the population at time  $t$  can be described as  $S(t) = \{X_{1,t}, X_{2,t}, \dots, X_{NP,t}\} = \{(H_{1,t}, \Lambda_{1,t}), (H_{2,t}, \Lambda_{2,t}), \dots, (H_{NP,t}, \Lambda_{NP,t})\}$ .

**Mutation.** The purpose of the mutation operation is to generate a new individual based on the different individuals.

According to formulas (12) and (13), we define the state

transition probability of mutation operation as follows:

$$P_{X_{i,t} \rightarrow V_{i,t}} = \alpha_i, \left( 0 \leq \alpha_i \leq \frac{1}{(NP-1) \cdot (NP-2) \cdot (NP-3)} \right), \quad (17)$$

**Crossover.** From formula (14), the probability of crossover operation is:

$$P_{(X_{i,t}, V_{i,t}) \rightarrow U_{i,t}} = k_i, (0 \leq k_i \leq 1). \quad (18)$$

**Selection.** From formula (16), the probability of selection operation generating evolved individual is:

$$P_{\text{select}} = \begin{cases} 1 & \text{if } \text{perf}(U_{i,t}) \geq \text{perf}(X_{i,t}) \\ 0 & \text{otherwise} \end{cases}. \quad (19)$$

Assuming that mutation, crossover, and selection operations of DE algorithm are independent of each other, and the selection process from three random individuals in mutation operation is also an independent random process. To sum up, the probability of individual transition from state  $X_{i,t}$  to state  $X_{i,t+1}$  in DE-Opt is:

$$P_{X_{i,t} \rightarrow X_{i,t+1}} = P_{X_{i,t} \rightarrow V_{i,t}} \cdot P_{(X_{i,t}, V_{i,t}) \rightarrow U_{i,t}} \cdot P_{\text{select}}. \quad (20)$$

From the above analyses, we note that for any two continuous population states set  $S(t) = \{X_{1,t}, X_{2,t}, \dots, X_{NP,t}\}$  and  $S(t+1) = \{X_{1,t+1}, X_{2,t+1}, \dots, X_{NP,t+1}\}$  in the population state sequence  $\{S(t) \mid t \in N^+\}$ , the state transition probability  $P(S(t) \rightarrow S(t+1))$  is only related to the population state at time  $t$ . It is proved that the population state sequence has Markov characteristics. Since the population state set is discrete and finite, and the probability in formulas (17), (18) and (19) are independent of  $t$ , the population state sequence  $\{S(t) \mid t \in N^+\}$  is a finite homogeneous Markov chain.

### 4.3.2 Convergence Criteria

For randomized algorithms, e.g., DE-Opt, the convergence criteria can be described as follows.

**Definition 1.** Given a new iteration solution  $S(t+1) = \Phi(S(t), \xi)$  of an optimization problem  $\langle I, \text{perf}(\cdot) \rangle$ , where  $\Phi(\cdot)$  is an optimization algorithm,  $\xi = \{S(1), S(2), \dots, S(t-1)\}$  is previous searched solutions,  $I$  is the feasible solution space. For DE-Opt,  $I$  is the hyperparameters searching space and  $\text{perf}(\cdot)$  is the fitness/performance function.

**Assumption 1.**  $\text{perf}(\Phi(S(t), \xi)) \geq \text{perf}(\xi)$ , s.t.  $\text{perf}(\Phi(S(t), \xi)) \geq \text{perf}(\xi)$ ,  $\xi \in I$ .

**Definition 2.** Population state space is  $S = \{S(t) = \{X_{i,t} \mid X_{i,t} \in I, 1 \leq i \leq NP\}\}$ , s.t.  $S(t) = \{X_{i,t} \mid X_{i,t} \in I, 1 \leq i \leq NP\}$  is the population state.

**Definition 3.** Global optimal state set is  $G = \{X \mid \text{perf}(X) = \text{perf}(g^*), X \in S\}$ , s.t.  $g^*$  is an optimal solution of  $\langle I, \text{perf}(\cdot) \rangle$ .

### 4.3.3 Convergence Proofs of DE-Opt

The DE-Opt can be represented as follows:

$$X_{i,t+1} = P_{X_{i,t} \rightarrow V_{i,t}} \cdot P_{(X_{i,t}, V_{i,t}) \rightarrow U_{i,t}} \cdot P_{\text{select}} \cdot X_{i,t}, \quad (21)$$

where  $i = \{1, 2, \dots, NP\}$ . Formula (21) can be proved to be convergent as follows.

**Lemma 1.** The DE-Opt satisfies assumption 1, i.e.  $\text{perf}(S(t+1)) \geq \text{perf}(S(t))$ .

**Proof.** According to selection operation, it is easy to proof

that  $\text{perf}(X_{i,t+1}) \geq \text{perf}(X_{i,t})$ . From definitions 1 and 2, it is also easy to proof that  $\text{perf}(S(t+1)) \geq \text{perf}(S(t))$ , i.e., DE-Opt is monotonically increasing, lemma 1 stands.

**Theorem 2.** The state sequence of DE-Opt  $\{S(t) \mid t \in N^+\}$  is a finite homogeneous irreducible aperiodic Markov chain.

**Proof.** From the Markov chain of DE-Opt,  $\{S(t) \mid t \in N^+\}$  has been proven to be a finite homogeneous Markov chain. From definition 2, we can obtain:

$$\sum S(t+1) \in S \cdot P\{S(t) \cdot P_{X_{i,t} \rightarrow X_{i,t+1}} = S(t+1)\} = 1, \quad (22)$$

where  $\forall S(t) \in S$ , therefore  $\{S(t) \mid t \in N^+\}$  is irreducible Markov chain. For the following formula:

$$P(S(t) \cdot P_{X_{i,t} \rightarrow X_{i,t+1}} = S(t+1)) = \sum S \cdot P_{X_{i,t} \rightarrow V_{i,t}} \cdot P_{(X_{i,t}, V_{i,t}) \rightarrow U_{i,t}} \cdot P_{\text{select}}, \quad (23)$$

there are

$$P_{X_{i,t} \rightarrow V_{i,t}} > 0, P_{(X_{i,t}, V_{i,t}) \rightarrow U_{i,t}} > 0, \text{ and } P_{\text{select}} > 0, \quad (24)$$

therefore

$$P(X_{i,t} \cdot P_{X_{i,t} \rightarrow X_{i,t+1}} = X_{i,t+1}) > 0. \quad (25)$$

Consequently,  $\{S(t) \mid t \in N^+\}$  is an aperiodic Markov chain.

**Theorem 3.** In DE-Opt,  $G$  from definition 3 is a closed set.

**Proof.** From the Chapman-Kolmogorov equation [31-32], transfer probability  $P_{S(a), S(b)}^l$  from state  $S(a)$  to state  $S(b)$  after  $l$  transfer steps is:

$$P_{S(a), S(b)}^l = \sum_{S(\varepsilon_1) \in S} \sum_{S(\varepsilon_2) \in S} \dots \sum_{S(\varepsilon_{l-1}) \in S} P(S(a) \cdot P_{X_{i,a} \rightarrow X_{i,\varepsilon_1}} = S(\varepsilon_1)) \cdot P(S(\varepsilon_1) \cdot P_{X_{i,\varepsilon_1} \rightarrow X_{i,\varepsilon_2}} = S(\varepsilon_2)) \dots P(S(\varepsilon_{l-1}) \cdot P_{X_{i,\varepsilon_{l-1}} \rightarrow X_{i,b}} = S(b)), \quad (26)$$

s.t.  $\forall l (l \geq 1), \forall S(a) = \{X_{1,a}, X_{2,a}, \dots, X_{NP,a}\} \in G, \forall S(b) \notin G$ .

In each step, there is a transfer from  $S(t-1) \in G$  to  $S(t)$ , and  $t \in \{1, 2, \dots, l\}$ . There is  $\text{perf}(S(t)) < \text{perf}(S(t-1)) = \text{perf}(g^*)$ , s.t.  $S(t-1) \in G$  and  $S(t) \notin G$ , therefore  $\exists P(S(t-1) \rightarrow S(t)) = 0$ , i.e.  $P_{S(a), S(b)}^l = 0$ . Consequently,  $G$  is a closed set on state space  $S$ .

**Theorem 4.** No non-empty closed set  $G'$  satisfies  $G' \cap G = \emptyset$  on population state space  $S$  of the DE-Opt.

**Proof.** (Proof by contradiction). Given a non-empty closed set  $G'$  on  $S$ ,  $G' \cap G = \emptyset$ ,  $S(a) = \{g_{1,a}, g_{2,a}, \dots, g_{NP,a}\} \in G$ ,  $\forall S(b) = \{X_{1,b}, X_{2,b}, \dots, X_{NP,b}\} \in G'$ , there is  $\text{perf}(X_{i,b}) < \text{perf}(g_{i,a})$ . Thus, from formula (24),  $\exists P(S(\varepsilon_{t+1}) \rightarrow S(\varepsilon_{t+a+1})) > 0$ ,  $P_{S(a), S(b)}^l > 0$ , which means that  $G'$  is not a non-empty closed set.

**Theorem 5** [31-32]. Transfer probability  $\pi_j$  satisfies  $\lim_{k \rightarrow \infty} p(S(k) = j) = \pi_j$  when  $j \in E$ , and  $\lim_{k \rightarrow \infty} p(S(k) = j) = 0$  when  $j \notin E$ , s.t.  $\exists$  non-empty closed sets  $E$  and  $V$  satisfy  $E \cap V = \emptyset$ .

**Theorem 6.** Let iterations tend to infinity,  $S$  will get into  $G$  in DE-Opt.

**Proof.** Theorems 2-5 ensure the authenticity of theorem 6.

Finally, from lemma 1 and theorems 2-6, it is easy to prove that DE-Opt is convergent in the given searching space. Therefore, Theorem 1 holds.

## 4.4 Algorithm Design and Complexity Analysis

Based on the analyses of Section IV.B, we design *Algorithm DE-Opt*, with its pseudo-code presented in Table 3. We also analyze its time complexity in improving a DL-based RS as follow. Supposing the time complexity of an original model with fixed  $\eta$  and  $\lambda$  is  $O(T_{\text{model}})$  and DE-Opt contains  $NP$  individuals. Then, the time complexity of DE-Opt is

$O(T_{model} \times NP)$ . Notably, an original model requires Grid Search to search optimal  $\eta$  and  $\lambda$ . Assuming that each hyperparameter needs  $\theta$  times search, then the final time complexity of an original model is  $O(T_{model} \times \theta^2)$ . In general,  $NP$  is set in a small range like from 5 to 20 (DE-Opt is 5, please refer to Section 5.5) [18-20], and each hyperparameter usually needs at least five times search (i.e.,  $\theta \geq 5$ ) [8, 22]. Hence, the time complexity of DE-Opt is less than that of Grid Search in searching  $\eta$  and  $\lambda$  in general. Besides, DE-Opt can auto-learn the optimal  $\eta$  and  $\lambda$ , which avoids the tedious Grid Search. In terms of space complexity, DE-Opt adopts  $NP$  arrays with length  $2 \times N$  to cache  $NP$  individuals, where each individual is a pair of  $\Lambda$  and  $H$ . Hence, the space complexity of DE-Opt is  $O(N \times NP)$ .

TABLE 3. ALGORITHM DE-OPT.

Steps	Operations
1	initializing
2	model construction
3	generate vectorized $H$ and $\Lambda$ as Population initialized for model
4	<b>while</b> epoch $t \leq \text{Epoch}_{\max}$ && not converge
5	calculate performance evolved in the <b>validation</b> set
6	<b>for</b> $i = 1$ to $NP$
7	initialize $X_i$ in $\text{Population}_{t-1}$
8	randomly select $X_{r1}$ , $X_{r2}$ and $X_{r3}$ in $\text{Population}_{t-1}$
9	compute $V_i$ by formula (12)
10	generate $U_i$ from $X_i$ and $V_i$ by formula (15)
11	$H_{\text{train}} = U_i^1$ , $\Lambda_{\text{train}} = U_i^2$
12	train model with $H_{\text{train}}$ and $\Lambda_{\text{train}}$ in the <b>training</b> set
13	calculate performance evolved in the <b>validation</b> set
14	Select $X_i$ by formula (16)
15	update $\text{Population}_{t-1}$ with $X_i$
16	<b>end for</b>
17	complete the evolution of $\text{Population}_{t-1}$ to $\text{Population}_t$
18	$t = t + 1$
19	<b>end while</b>

## 5 EXPERIMENTS

In the subsequent experiments, we aim at answering the following research questions (RQs):

- RQ.1. Can DE-Opt boost state-of-the-art recommendation models with fixed hyperparameters in both rating prediction and item ranking tasks?
- RQ.2. Does DE-Opt significantly outperform state-of-the-art hyperparameter searching competitors?
- RQ.3. How does DE-Opt empirically converge to optimal hyperparameter sets?
- RQ.4. How does the individual number  $NP$  of DE influence the performance of DE-Opt?
- RQ.5. How do mutation and crossover strategies of DE influence the performance of DE-Opt?

### 5.1 General Settings

**Datasets.** We adopt three real-world datasets [5] of *FilmTrust*<sup>1</sup>, *MovieLens 1M*<sup>1</sup>, and *Each Movie*<sup>1</sup> as the benchmark in the experiments. Their statistics are summarized in Table 4. We adopt the explicit ratings for the rating prediction task and transform them into binary code for the item ranking task.

TABLE 4. STATISTICS OF THE STUDIED DATASETS.

No.	Name	#User	#Item	#Interaction	Density*
D1	Film Trust	1,508	2,071	37,919	1.21%
D2	MovieLens 1M	6,040	3,706	1,000,209	4.47%
D3	Each Movie	72,916	1,628	2,811,983	2.36%

\* Density denotes the percentage of observed entries in the dataset.

**Evaluation Metrics.** Missing rating prediction and item ranking prediction (top- $K$  recommendation) are two common tasks in an RS [5]. Considering missing rating prediction, mean absolute error (MAE) and root mean squared error (RMSE) are widely adopted as the evaluation metrics [5, 22] as follows:

$$MAE = \left( \sum_{(u,j) \in \Gamma} |y_{u,j} - \hat{y}_{u,j}|_{abs} \right) / |\Gamma|,$$

$$RMSE = \sqrt{\left( \sum_{(u,j) \in \Gamma} (y_{u,j} - \hat{y}_{u,j})^2 \right) / |\Gamma|},$$

where  $\Gamma$  denotes the testing set and  $|\cdot|_{abs}$  calculates the absolute value of a given number. Note that a lower MAE/RMSE indicates a higher rating prediction accuracy.

On the other hand, item ranking prediction produces a ranked list with  $K$  items for recommendation, where  $K$  is a cutoff hyperparameter. Normalized discounted cumulative gain (NDCG), Recall, and hit ratio (HR) are widely adopted to evaluate the ranking prediction accuracy [8, 14]. Let  $\Gamma_K(u)$  denote the intersection set between the top- $K$  items generated by the tested model and the testing items of  $\Gamma$  on user  $u$ . Then, the NDCG of a tested model is computed as follows:

$$DCG@K(u, \hat{Y}) = \sum_{j=1}^K \frac{2^{\hat{y}_{u,j}} - 1}{\log_2(j+1)},$$

$$DCG@K(u, \Gamma) = \sum_{j=1}^K \frac{2^{y_{u,j}} - 1}{\log_2(j+1)},$$

$$NDCG@K = \frac{1}{|\Psi|} \sum_{u \in \Psi} \left( \frac{DCG@K(u, \hat{Y})}{DCG@K(u, \Gamma)} \right),$$

where  $\hat{y}_{u,j}$  denotes the  $j$ -th prediction in  $\Gamma_K(u)$  in descending order,  $y_{u,j}$  denotes the actual  $j$ -th item of  $\Gamma$  on user  $u$  in descending order, and  $\Psi$  denotes a set of contained users of  $\Gamma$ . The Recall of a tested model is computed as follows:

$$Recall@K = \frac{1}{|\Psi|} \sum_{u \in \Psi} \frac{|\Gamma_K(u)|}{K}.$$

If a test item appears in the recommended list, it is deemed as a hit. Then, the HR is calculated as follows:

$$hits@K(u) = \begin{cases} 1 & \Gamma_K(u) \not\subseteq \emptyset \\ 0 & \Gamma_K(u) \subseteq \emptyset \end{cases}$$

$$HR@K = \frac{1}{|\Psi|} \sum_{u \in \Psi} hits@K(u).$$

Note that NDCG, Recall, and HR lie in the scale of  $[0, 1]$ , where a higher value indicates a higher ranking prediction accuracy.

**Baselines.** We compare our DE-Opt with fourteen state-of-the-art models, including six fixed hyperparameter models and eight adaptive hyperparameter models. Their details are summarized in Table 5.

TABLE 5. SUMMARY OF COMPARISON MODELS.

Model	Description
AutoRec [22]	A classic DL-based rating prediction model. It is a representative collaborative filtering method based on autoencoder. (WWW 2015)
NRT [23]	It is a multitask learning with a rating prediction framework developed by combining multilayer perceptron and recurrent neural networks. (SIGIR 2017)
Fixed $\eta$ and $\lambda$ model	
MetaMF [24]	It is a federated learning framework for private rating prediction for mobile environments. (SIGIR 2020)
NeuMF [14]	A classic DL-based item ranking model. It combines traditional matrix factorization and multilayer perceptron. (WWW 2017)
LRML [25]	It is a neural architecture for collaborative ranking via memory-based attention. (WWW 2018)
NGCF [21]	A GCN-based item ranking model. It exploits user-item graph structure by propagating. (SIGIR 2019)
$\eta$ -adaptive model	
Adam [17]	A classic SGD-based $\eta$ -adaptive optimizer. It only requires first-order gradients with little memory for efficient stochastic optimization. (ICLR 2015)
AMSGrad [12]	It is a variant of Adam that fixes the convergence issue and improves empirical performance. (ICLR 2018)
AdaMod [13]	It proposes automatic warmup heuristic and long-term learning rate buffer to improve Adam. (arXiv 2019)
$\lambda$ -adaptive model	
SGDA [11]	It makes $\lambda$ adaptive based on dimension-wise regularization and SGD. It is limited to adopt $\eta$ -adaptive optimizers [8]. (WSDM 2012)
$\lambda$ Opt [8]	It updates $\lambda$ automatically and adaptively during training. It allows fine-grained regularization to build a better generalized model. (KDD 2019)
$\eta$ - $\lambda$ -adaptive model	
$\lambda$ Opt-Adam	We adopt $\eta$ -adaptive Adam to optimize $\lambda$ Opt.
$\lambda$ Opt-AMSGrad	We adopt $\eta$ -adaptive AMSGrad to optimize $\lambda$ Opt.
$\lambda$ Opt-AdaMod	We adopt $\eta$ -adaptive AdaMod to optimize $\lambda$ Opt.

**Implementation Detail.** We set Grid Search and our DE-Opt to have the same searching space. Specially, we adopted the ten times amplification principle to set their searching space, i.e., the optimal values searched by Grid Search are amplified ten times for setting DE-Opt's searching space. Please refer to Table S.2-4 in the Supplementary File to see the details of the searching space. We randomly split the observed entries of each dataset into the training-validation-testing sets with 70%-10%-20% ratio, respectively. DE-Opt auto-learns the DL-based RSs' hyperparameters on the training set and the learning performances are supervised on the validation set, i.e., the DL-based RSs' hyperparameters are auto-learned on the training set to make the resultant RS perform best on the validation set. Finally, all the trained models are tested on the testing set. All the other hyperparameters of tested models (except for  $\lambda$  and  $\eta$ ) are set by the original papers' recommendations. SGD is the default optimizer. Adaptive hyperparameter models have a basic deep matrix factorization module. The prediction rating and ranking list for a particular item are generated according to the original papers. We repeat each experiment five times and report average results.

**System Configuration.** The experiments are implemented by Python 3.7.0 with the libraries of Numpy, Tensorflow, Scipy, Sklearn, Pandas, and Torch. All experiments are run by 11G GTX1080ti GPUs. Our resource code is available at the following link: <https://github.com/Wuziqiao/DE-Opt>.

## 5.2 Comparison with Fixed Hyperparameter Model (RQ. 1)

### 5.2.1 Rating Prediction Task

Table 6 presents the comparison results of rating prediction, where we observe that DE-Opt can achieve lower RMSE/MAE than Grid Search method except for only one case. To check whether DE-Opt significantly outperforms Grid Search method, we conduct the Wilcoxon signed-ranks test [26] on the comparisons of RMSE/MAE of Table 6. The  $p$ -value denotes the significance level, and if it is lower than a significance level 0.05,

TABLE 6. THE COMPARISON RESULT OF RATING PREDICTION ACCURACY.

Dataset	Metric	AutoRec		NRT		MetaMF	
		Grid Search	DE-Opt	Grid Search	DE-Opt	Grid Search	DE-Opt
D1	RMSE	0.8223	0.8225	0.8087	<b>0.8023</b>	0.8239	<b>0.8214</b>
	MAE	0.6188	<b>0.6179</b>	0.6259	<b>0.6124</b>	0.6469	<b>0.6379</b>
D2	RMSE	0.8532	<b>0.8520</b>	0.8791	<b>0.8771</b>	0.9370	<b>0.9090</b>
	MAE	0.6665	<b>0.6659</b>	0.6919	<b>0.6893</b>	0.7433	<b>0.7200</b>
D3	RMSE	0.3170	<b>0.3102</b>	<b>0.2594</b>	0.2599	0.2815	<b>0.2771</b>
	MAE	0.2370	<b>0.2359</b>	0.2021	<b>0.2018</b>	0.2240	<b>0.2197</b>
Statistics		5/1		5/1		6/0	
		<b>0.0313</b>		<b>0.0469</b>		<b>0.0156</b>	

\*A The accepted hypotheses with a significance level of 0.05 are highlighted.

TABLE 7. THE COMPARISON RESULT OF ITEM RANKING ACCURACY.

Dataset	Metric	NeuMF		LRML		NGCF	
		Grid Search	DE-Opt	Grid Search	DE-Opt	Grid Search	DE-Opt
D1	NDCG@10	0.4513	<b>0.4641</b>	<b>0.3781</b>	0.3746	0.3734	<b>0.3773</b>
	NDCG@20	0.5377	<b>0.5464</b>	<b>0.4588</b>	0.4556	0.4562	<b>0.4602</b>
	Recall@10	0.5404	<b>0.5500</b>	0.4385	<b>0.4460</b>	<b>0.4002</b>	0.3987
	Recall@20	0.7547	<b>0.7587</b>	0.6539	<b>0.6619</b>	<b>0.5778</b>	0.5775
D2	NDCG@10	0.3035	<b>0.3047</b>	0.2136	<b>0.2234</b>	0.3541	<b>0.3544</b>
	NDCG@20	0.2953	<b>0.2990</b>	0.2098	<b>0.2308</b>	0.3852	<b>0.3857</b>
	Recall@10	0.1241	<b>0.1262</b>	0.0806	<b>0.1296</b>	0.1366	<b>0.1392</b>
	Recall@20	0.1972	<b>0.2026</b>	0.1350	<b>0.2068</b>	0.2167	<b>0.2211</b>
D3	NDCG@10	0.1078	<b>0.1105</b>	0.0655	<b>0.0752</b>	<b>0.0523</b>	0.0528
	NDCG@20	0.1460	<b>0.1481</b>	0.0921	<b>0.1059</b>	0.0747	<b>0.0753</b>
	Recall@10	0.1908	<b>0.1921</b>	0.1138	<b>0.1315</b>	0.0651	<b>0.0663</b>
	Recall@20	<b>0.3158</b>	0.3147	0.2010	<b>0.2321</b>	0.1164	<b>0.1180</b>
Statistics		11/1		10/2		9/3	
		<b>0.0005</b>		<b>0.0012</b>		<b>0.0081</b>	

\*A The accepted hypotheses with a significance level of 0.05 are highlighted.

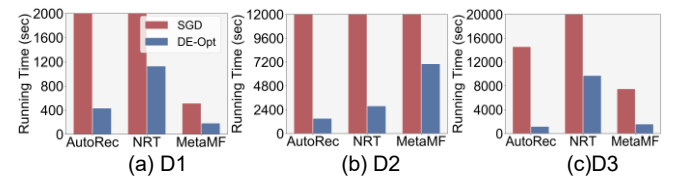


Fig. 4. The comparison results of GPU running time in rating prediction task.

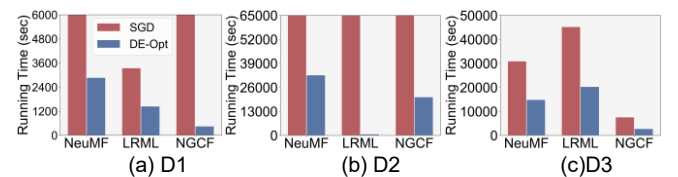


Fig. 5. The comparison results of GPU running time in item ranking task.



the hypotheses are accepted. The results demonstrate that DE-Opt can significantly improve AutoRec, NRT, and MetaMF in rating prediction accuracy.

Besides, Figure 4 records the GPU running time of each model tested on all the datasets. We observe that Grid Search significantly consumes more running time than DE-Opt. The reason is that DE-Opt can auto-learn the appropriate hyperparameters with only one-time training, which avoids the multiple times training with Grid Search.

### 5.2.2 Item Ranking Task

The accuracy comparison results of item ranking are shown in Table 7. They are also analyzed with the Wilcoxon signed-ranks test [26] with a significance level 0.05. Similar to rating prediction results, the results validate that DE-Opt can significantly boost NeuMF, LRML, and NGCF in item ranking. Besides, the comparison results of GPU running time of item ranking are also shown in Figure 5. The results are also similar to rating prediction results, i.e., DE-Opt consumes significantly less GPU running time than Grid Search.

### 5.3 Comparison with Adaptive Hyperparameter Model (RQ. 2)

The accuracy comparison results of item ranking are shown in Table 8. We also adopt the Wilcoxon signed-ranks test [26] to check these comparison results with a significance level 0.05. Besides, we further make the Friedman test [26]. It is to simultaneously compare the performance of multiple models on multiple datasets by the F-rank value. From Table 8, we have the important findings:

- DE-Opt achieves the higher ranking prediction accuracy on most cases. It loses the comparison with only 3 cases. Specifically, the total win/ loss cases situation is 92/4.
- All the  $p$ -values are smaller than 0.05, which indicates that DE-Opt has significantly higher ranking prediction accuracy than all the adaptive hyperparameter models with a

significance level 0.05.

- DE-Opt has the highest F-rank value among all the models, which further verifies that it has the overall highest ranking prediction accuracy on the tested datasets.

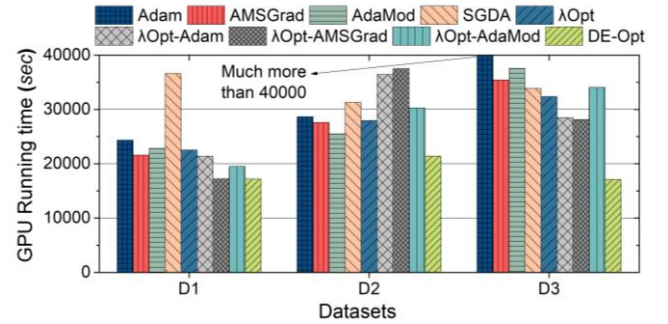


Fig. 6. The comparison results of GPU running time between DE-Opt and adaptive hyperparameter models.

Besides, the computational efficiency is also compared. Figure 6 shows the GPU running time of all the involved models, where we note that DE-Opt consumes significantly less running time than the other adaptive hyperparameter models.

Therefore, the above experimental results demonstrate that DE-Opt significantly outperforms the adaptive hyperparameter models in terms of both prediction accuracy and computational efficiency.

### 5.4 Auto-learning the Hyperparameters (RQ. 3)

To check how DE-Opt auto-learns  $\lambda$  and  $\eta$ , we draw their change curves during its training process. A deep matrix factorization is also adopted as the basic model of DE-Opt. Figures 7 and 8 show the results on all the datasets, where we find that both  $\lambda$  and  $\eta$  oscillate at the beginning of training. Since DE-Opt tries to search for better hyperparameters in the initial steps, the oscillation can prompt hyperparameters to jump out of the poor local optimal solutions. As the training epoch increases, they gradually stabilize in a certain range. The

TABLE 8. THE COMPARISON RESULTS WITH ADAPTIVE HYPERPARAMETER MODELS.

Dataset	Metrics	$\eta$ -adaptive( $\lambda$ -GridSearch)			$\lambda$ -adaptive( $\eta$ -GridSearch)		$\eta$ - $\lambda$ -adaptive			DE-Opt
		Adam	AMSGrad	AdaMod	SGDA	$\lambda$ Opt	$\lambda$ Opt-Adam	$\lambda$ Opt-AMSGrad	$\lambda$ Opt-Ada-Mod	
D1	NDCG@10	0.4321•	0.4412•	0.4311•	0.4115•	<b>0.4466</b>	0.4332•	0.4401•	0.4342•	0.4438
	NDCG@20	0.4924•	0.4809•	0.4723•	0.4660•	0.3955•	0.4651•	0.4863•	0.4628•	<b>0.5061</b>
	HR@10	<b>0.6446</b>	0.6370•	0.6324•	0.5907•	0.5995•	0.6329•	0.6407•	0.6350•	0.6419
	HR@20	0.7920•	0.7891•	0.7787•	0.7404•	0.6800•	0.7778•	0.7840•	0.7814•	<b>0.8044</b>
D2	NDCG@10	0.0482•	0.0474•	0.0476•	0.0031•	0.0238•	0.0852•	<b>0.1095</b>	0.0863•	0.1050
	NDCG@20	0.0627•	0.0560•	0.0642•	0.0074•	0.0323•	0.1033•	0.1318•	0.1048•	<b>0.1340</b>
	HR@10	0.0965•	0.0951•	0.0970•	0.0069•	0.0482•	0.1594•	<b>0.1959</b>	0.1611•	0.1907
	HR@20	0.1579•	0.1421•	0.1595•	0.0225•	0.0831•	0.2370•	0.2826•	0.2402•	<b>0.2847</b>
D3	NDCG@10	0.0697•	0.0557•	0.0689•	0.0297•	0.0222•	0.0614•	0.0659•	0.0621•	<b>0.0778</b>
	NDCG@20	0.0967•	0.0973•	0.0987•	0.0417•	0.0362•	0.0859•	0.0943•	0.0913•	<b>0.1013</b>
	HR@10	0.1478•	0.1183•	0.1474•	0.0624•	0.0495•	0.1309•	0.1407•	0.1315•	<b>0.1640</b>
	HR@20	0.2593•	0.2579•	0.2609•	0.1097•	0.1059•	0.2327•	0.2546•	0.2453•	<b>0.2701</b>
Statistics	Win/Loss	11/1	12/0	12/0	12/0	11/1	12/0	10/2	12/0	92/4○
	$p$ -value*	<b>0.0005</b>	<b>0.0002</b>	<b>0.0002</b>	<b>0.0002</b>	<b>0.0005</b>	<b>0.0002</b>	<b>0.0134</b>	<b>0.0002</b>	—
	F-rank◇	6.17	4.75	5.25	1.67	2.08	4.33	6.83	5.25	<b>8.67</b>

• Indicates DE-Opt outperforms the comparison models. ○ The total Win/Loss cases of DE-Opt. \* The accepted hypotheses with a significance level of 0.05 are highlighted. ◇ A higher F-rank value denotes a higher item ranking accuracy.

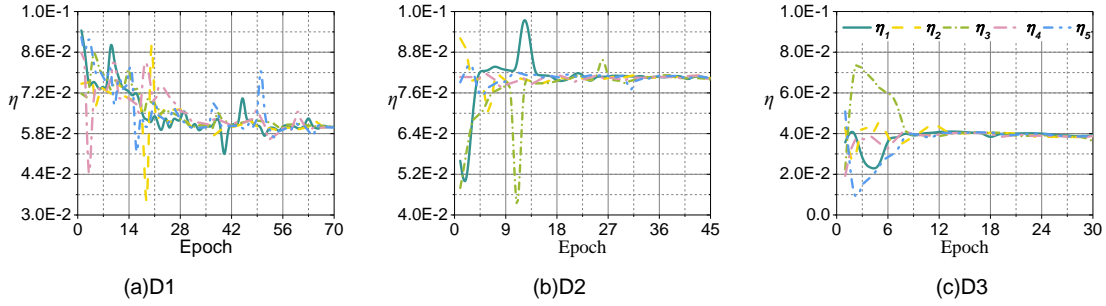


Fig. 7. Auto-learning trajectories of  $\eta$  with five individuals on all the datasets.

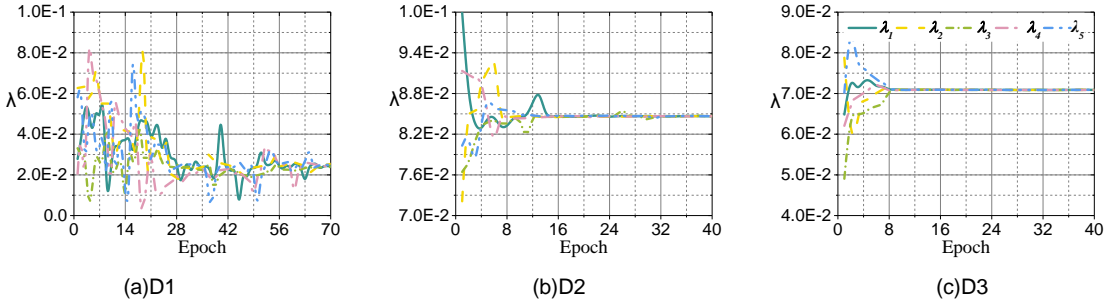


Fig. 8. Auto-learning trajectories of  $\lambda$  with five individuals on all the dataset.

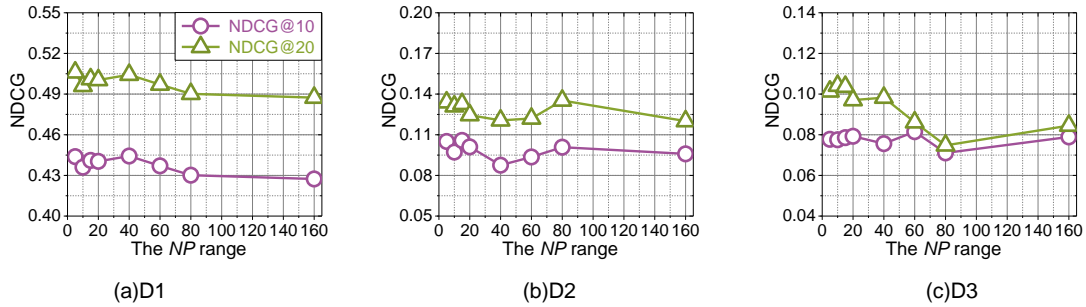


Fig. 9. The influence of  $NP$  on NDCG on all the dataset.

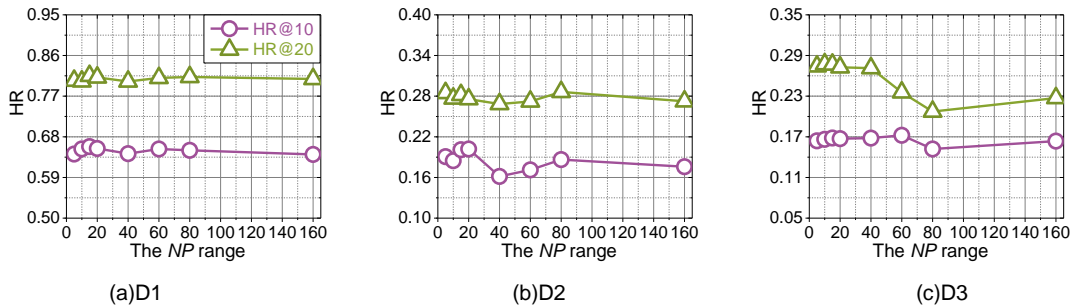


Fig. 10. The influence of  $NP$  on HR on all the dataset.

above results well prove that DE-Opt can make a DL-based recommendation model'  $\lambda$  and  $\eta$  be simultaneously auto-learned at layer-granularity during its training process. Besides, we also check the descent/ascent curves of loss/accuracy of the tested model during the training process. Please refer to Figures S.1–3 in the Supplementary File to see the details. The results also prove that DE-Opt can make the loss/accuracy of the tested model better with more training epochs until convergence.

### 5.5 Influence of Individual Number $NP$ (RQ. 4)

To evaluate the influences of individual number  $NP$  to DE-Opt, we carry out experiments with different  $NP$  settings. The results are presented in Figures 9 and 10. We note that DE-Opt performs relatively robustly and well when  $NP$  is between 5 to

15. Subsequently, as  $NP$  continues to increase, the NDCG and HR decreases in general. The reason may be that more individuals (larger  $NP$ ) cause more random processes during the training process. Overall, DE-Opt is relatively robust to the individual number  $NP$ . We set  $NP=5$  for our DE-Opt.

### 5.6 Influences of Mutation and Crossover Strategies (RQ. 5)

This sub-section analyzes the influences of different mutation and crossover strategies to DE-Opt. Concretely, we compare the six classical mutation strategies [42–45], i.e., DE/Rand/1 (adopted in Section IV.B), DE/Best/1, DE/RandToBest/1, DE/Best/2, DE/Rand/2, and DE/RandToBest/2; and the three situations of crossover strategies [42–45], i.e., without crossover,

TABLE 9. THE CMPARISON RESULT WITH SX MTATION STRATEGIES IN NDCG AND HR.

Dataset	Metric	DE/Rand/1	DE/Best/1	DE/RandToBest/1	DE/Best/2	DE/Rand/2	DE/RandToBest/2
D1	NDCG@10	0.4438•	0.4504•	<b>0.4777</b>	0.4538•	0.4439•	0.4613
	NDCG@20	0.5061•	0.5057•	0.5021•	0.5044•	0.5076•	<b>0.5170</b>
	HR@10	0.6419•	0.6450•	0.6503•	0.6511•	0.6461•	<b>0.6615</b>
	HR@20	0.8044•	0.8084•	0.8127•	0.8140•	0.7988•	<b>0.8228</b>
D2	NDCG@10	0.1050•	0.1034•	<b>0.1098</b>	0.1097	0.0956•	0.1053
	NDCG@20	0.1340•	0.1375•	0.1387•	0.1353•	0.1285•	<b>0.1390</b>
	HR@10	0.1907•	0.1917•	0.2010	<b>0.2028</b>	0.1916•	0.1998
	HR@20	0.2847•	<b>0.2945</b>	0.2838•	0.2857•	0.2825•	0.2858
D3	NDCG@10	0.0778•	0.0785•	0.0753•	0.0759	0.0728•	<b>0.0786</b>
	NDCG@20	0.1013•	0.1032•	<b>0.1088</b>	0.1078•	0.1060•	0.1078
	HR@10	0.1640•	0.1640•	0.1612•	0.1634•	0.1530•	<b>0.1705</b>
	HR@20	0.2701•	0.2724•	0.2794•	0.2730•	0.2658•	<b>0.2819</b>
Statistics	Win/Loss	12/0	11/1	8/4	9/3	12/0	52/8○
	p-value*	<b>0.0002</b>	<b>0.0046</b>	0.1331	<b>0.0147</b>	<b>0.0002</b>	—
	F-rank◇	2.375	3.375	4.0	4.042	1.83	5.375

binomial crossover, and arithmetic crossover (adopted in Section IV.B). Next, we mainly analyze how these strategies influence the prediction accuracy of DE-Opt. Note that the influence of computational efficiency is not significant, we move the results to Table S.5 in the Supplementary File. Besides, the details of these mutation and crossover strategies are also moved to Section S.IV in the Supplementary File.

### 5.6.1 Influences of Mutation Strategies

Table 9 records the comparison results of different mutation strategies on all the datasets, where we observe that DE/RandToBest/2 seems to perform best. To further validate this observation, we conduct the statistics of win/loss counts, the Wilcoxon signed-ranks test, and the Friedman test by comparing DE/RandToBest/2 with the other five mutation strategies. The statistical results are also recorded in Table 9, where we note that: 1) DE/RandToBest/2 evidently wins the other mutation strategies on most cases, 2) DE/RandToBest/2 significantly outperforms the other mutation strategies except for DE/RandToBest/1, and 3) DE/RandToBest/2 performs best among all the mutation strategies as it achieves the highest F-rank value.

### 5.6.2 Influences of Crossover Strategies

Table 10 presents the comparison results of different crossover strategies on all the datasets. We also conduct the statistics of win/loss counts, the Wilcoxon signed-ranks test, and the Friedman test by comparing arithmetic crossover with the other two strategies. The statistical results evidently show that arithmetic crossover significantly outperforms the other two strategies as the  $p$ -values are smaller than 0.05, and it achieves the highest F-rank value.

### 5.6.3 Summary

By analyzing the different mutation and crossover strategies, we empirically demonstrate that DE/RandToBest/2 mutation

TABLE 10. THE CMPARISON RESULT OF IEM RNKING ACURACY.

Dataset	Metric	WithOut Cross-over	Binomial Crossover	Arithmetic Crossover
D1	NDCG@10	0.3898•	0.4067•	<b>0.4438</b>
	NDCG@20	0.4386•	0.4661•	<b>0.5061</b>
	Recall@10	0.5603•	0.6190•	<b>0.6419</b>
	Recall@20	0.7325•	0.7647•	<b>0.8044</b>
D2	NDCG@10	0.0917•	<b>0.1080</b>	0.1050
	NDCG@20	0.1312•	0.1189•	<b>0.1340</b>
	Recall@10	0.1706•	0.1896•	<b>0.1907</b>
	Recall@20	0.2829•	0.2596•	<b>0.2847</b>
D3	NDCG@10	0.0461•	0.0712•	<b>0.0778</b>
	NDCG@20	0.0924•	<b>0.1041</b>	0.1013
	Recall@10	0.1019•	0.1538•	<b>0.1640</b>
	Recall@20	0.2570•	0.2618•	<b>0.2701</b>
Statistics	Win/Loss	12/0	10/2	22/2○
	p-value*	<b>0.0002</b>	<b>0.0024</b>	—
	F-rank◇	1.177	2	2.833

and arithmetic crossover strategies make DE-Opt perform better than the other strategies. Therefore, we recommend adopting DE/RandToBest/2 mutation and arithmetic crossover strategies for DE-Opt in practical applications.

## 6 CONCLUSION

This paper proposes a general hyperparameter optimization framework based on differential evolution (DE) for existing deep learning (DL)-based recommender systems (RSs), named DE-Opt. Its main idea is to incorporate DE into a DL-based RS's training process to auto-learn its hyperparameters  $\lambda$  (regularization coefficient) and  $\eta$  (learning rate) at once. DE-Opt allows to yield optimal hyperparameter pairs at a layer-granularity, where each layer has different  $\lambda$  and  $\eta$ , thereby uplifting the resultant recommendation performance. The empirical evaluation of DE-Opt was benchmarked on three

widely-used datasets with extensive experiments. The results substantiate that: 1) DE-Opt is compatible with six state-of-the-art DL-based recommendation models and significantly improves their performance while making their  $\lambda$  and  $\eta$  adaptive, and 2) DE-Opt significantly outperforms its competitors whose  $\lambda$  and/or  $\eta$  are adaptive. In the future, we will extend DE-Opt to more advanced RS models to evaluate its performance in complex scenarios, such as agricultural RS [47].

## ACKNOWLEDGMENT

This work is supported in part by the National Natural Science Foundation of China under grants 62072429 and 62293501, and in part by the Key Cooperation Project of Chongqing Municipal Education Commission (HZ2021008, HZ2021017).

## REFERENCES

- [1] Y. Gu, K. Vyas, M. Shen, J. Yang, and G. Yang, "Deep graph-based multimodal feature embedding for endomicroscopy image retrieval," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 32, no. 2, pp. 481–492, 2021.
- [2] L. Wu, X. He, X. Wang, K. Zhang and M. Wang, "A Survey on Accuracy-oriented Neural Recommendation: From Collaborative Filtering to Information-rich Recommendation," *IEEE Trans. Knowl. Data Eng.*, pp. 1–1, 2022. doi: 10.1109/TKDE.2022.3145690.
- [3] D. Wang, X. Zhang, D. Yu, G. Xu, and S. Deng, "CAME: content- and context-aware music embedding for recommendation," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 32, no. 3, pp. 1375–1388, 2021.
- [4] L. Ren and W. Wang, "A Granular SVM-Based Method for Top-N Web Services Recommendation," *IEEE Trans. Services Comput.*, vol. 15, no. 1, pp. 457–469, 2022.
- [5] S. Zhang, L. Yao, A. Sun, and Y. Tay, "Deep learning based recommender system: A survey and new perspectives," *ACM Comput. Surv.*, vol. 52, no. 1, pp. 5:1–5:38, 2019.
- [6] L. Huang, Y. Ma, S. Wang and Y. Liu, "An Attention-Based Spatiotemporal LSTM Network for Next POI Recommendation," *IEEE Trans. Services Comput.*, vol. 14, no. 6, pp. 1585–1597, 2021.
- [7] M. Maher and S. Sakr, "Smartml: A meta learning-based framework for automated selection and hyperparameter tuning for machine learning algorithms," in *Proceedings of the 22nd International Conference on Extending Database Technology, EDBT*, 2019, pp. 554–557.
- [8] Y. Chen, B. Chen, X. He, C. Gao, Y. Li, J. Lou, and Y. Wang, "llopt: Learn to regularize recommender models in finer levels," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA*, 2019, pp. 978–986.
- [9] X. Dong, J. Shen, W. Wang, L. Shao, H. Ling, and F. Porikli, "Dynamical hyperparameter optimization via deep reinforcement learning in tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 5, pp. 1515–1529, 2021.
- [10] E. Ndiaye, T. Le, O. Fercoq, J. Salmon, and I. Takeuchi, "Safe grid search with optimal complexity," in *Proceedings of the 36th International Conference on Machine Learning, ICML, Long Beach, California, USA*, 2019, pp. 4771–4780.
- [11] S. Rendle, "Learning recommender systems with adaptive regularization," in *Proceedings of the 15th International Conference on Web Search and Web Data Mining, WSDM*, 2012, pp. 133–142.
- [12] S. J. Reddi, S. Kale, and S. Kumar, "On the convergence of adam and beyond," in *Proceedings of the 6th International Conference on Learning Representations, ICLR, Vancouver, BC, Canada*, 2018.
- [13] J. Ding, X. Ren, R. Luo, and X. Sun, "An adaptive and momental bound method for stochastic learning," *CoRR*, vol. abs/1910.12249, 2019. Available: <http://arxiv.org/abs/1910.12249>
- [14] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T. Chua, "Neural collaborative filtering," in *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia*, 2017, pp. 173–182.
- [15] Z. Fu, H. Gao, W. Guo, S. K. Jha, J. Jia, X. Liu, B. Long, J. Shi, S. Wang, and M. Zhou, "Deep learning for search and recommender systems in practice," in *Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA*, 2020, pp. 3515–3516.
- [16] H. Lee, J. Im, S. Jang, H. Cho, and S. Chung, "Melu: Meta-learned user preference estimator for cold-start recommendation," in *Proceedings of the 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Anchorage, AK, USA*, 2019, pp. 1073–1082.
- [17] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA*, 2015.
- [18] Bilal, M. Pant, H. Zaheer, L. García-Hernández, and A. Abraham, "Differential evolution: A review of more than two decades of research," *Eng. Appl. Artif. Intell.*, vol. 90, p. 103479, 2020.
- [19] I. Triguero, S. García, and F. Herrera, "SEG-SSC: A framework based on synthetic examples generation for self-labeled semi-supervised classification," *IEEE Trans. Cybern.*, vol. 45, no. 4, pp. 622–634, 2015.
- [20] F. Neri and V. Tirronen, "Scale factor local search in differential evolution," *Memetic Comput.*, vol. 1, no. 2, pp. 153–171, 2009.
- [21] X. Wang, X. He, M. Wang, F. Feng, and T. Chua, "Neural graph collaborative filtering," in *Proceedings of the 42nd International Conference on Research and Development in Information Retrieval, ACM SIGIR, Paris, France*, July 21–25, 2019, pp. 165–174.
- [22] S. Sedhain, A. K. Menon, S. Sanner, and L. Xie, "Autorec: Autoencoders meet collaborative filtering," in *Proceedings of the 24th International Conference on World Wide Web Companion, WWW 2015, Florence, Italy, May 18–22*, 2015, pp. 111–112.
- [23] P. Li, Z. Wang, Z. Ren, L. Bing, and W. Lam, "Neural rating regression with abstractive tips generation for recommendation," in *Proceedings of the 40th International Conference on Research and Development in Information Retrieval, ACM SIGIR, Shinjuku, Tokyo, Japan, August 7–11*, 2017, pp. 345–354.
- [24] Y. Lin, P. Ren, Z. Chen, Z. Ren, D. Yu, J. Ma, M. de Rijke, and X. Cheng, "Meta matrix factorization for federated rating predictions," in *Proceedings of the 43rd International conference on research and development in Information Retrieval, ACM SIGIR, Virtual Event, China, July 25–30*, 2020, pp. 981–990.
- [25] Y. Tay, L. A. Tuan, and S. C. Hui, "Latent relational metric learning via memory-based attention for collaborative ranking," in *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW, Lyon, France, April 23–27*, 2018, pp. 729–739.
- [26] J. Demsar, "Statistical comparisons of classifiers over multiple data sets," *J. Mach. Learn. Res.*, vol. 7, pp. 1–30, 2006.
- [27] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, "Lightgcn: Simplifying and powering graph convolution network for recommendation," in *Proceedings of the 43rd International Conference on Research and Development in Information Retrieval, ACM SIGIR, Virtual Event, China, July 25–30*, 2020, pp. 639–648.
- [28] X. Wang, R. Wang, C. Shi, G. Song, and Q. Li, "Multi-component graph convolutional collaborative filtering," in *Proceedings of the 34th AAAI Conference on Artificial Intelligence, AAAI, New York, NY, USA, February 7–12*, 2020, pp. 6267–6274.
- [29] G. Xu and G. Yu, "On convergence analysis of particle swarm optimization algorithm," *J. Comput. Appl. Math.*, vol. 333, pp. 65–73, 2018.
- [30] J. Kwon, "Robust visual tracking based on variational auto-encoding markov chain monte carlo," *Inf. Sci.*, vol. 512, pp. 1308–1323, 2020.
- [31] Q. Li and M. Shang, "BALFA: A brain storm optimization-based adaptive latent factor analysis model," *Inf. Sci.*, vol. 578, pp. 913–929, 2021.
- [32] R. Poli and W. B. Langdon, "Markov chain models of bare-bones particle swarm optimizers," *Genetic and Evolutionary Computation Conference, GECCO 2007, Proceedings, London, England, UK, July 7–11*, 2007, pp. 103–110.



ACM, 2007, pp. 142–149.

- [33] S. Roy, S. Sridharan, S. Jain, and A. Raghunathan, "Txsim: Modeling training of deep neural networks on resistive crossbar systems," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 29, no. 4, pp. 730–738, 2021.
- [34] M. Wang, W. Fu, X. He, S. Hao, and X. Wu, "A survey on large-scale machine learning," *IEEE Trans. Knowl. Data Eng.*, pp. 1–1, 2020. DOI:10.1109/TKDE.2020.3015777
- [35] R. Yu, Y. Gong, X. He, Y. Zhu, Q. Liu, W. Ou, and B. An, "Personalized adaptive meta learning for cold-start user preference prediction," in *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, AAAI, Virtual Event, February 2-9, 2021. pp. 10 772–10 780.
- [36] D. Wu, Q. He, X. Luo, M. Shang, Y. He and G. Wang, "A Posterior-Neighborhood-Regularized Latent Factor Model for Highly Accurate Web Service QoS Prediction," *IEEE Trans. Services Comput.*, vol. 15, no. 2, pp. 793–805, 2022.
- [37] D. Wu, M. Shang, X. Luo, and Z. Wang, "An L1-and-L2-norm-oriented latent factor model for recommender systems," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 33, no. 10, pp. 5775–5788, 2022.
- [38] R. P. Parouha and P. Verma, "State-of-the-art reviews of meta-heuristic algorithms with their novel proposal for unconstrained optimization and applications," *Archives of Computational Methods in Engineering*, pp. 1–67, 2021.
- [39] D. Gong, Z. Zhang, Q. Shi, A. van den Hengel, C. Shen, and Y. Zhang, "Learning deep gradient descent optimization for image deconvolution," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 31, no. 12, pp. 5468–5482, 2020.
- [40] Q. Zhang, L. Cao, C. Shi and Z. Niu, "Neural Time-Aware Sequential Recommendation by Jointly Modeling Preference Dynamics and Explicit Feature Couplings," *IEEE Trans. Neural Networks Learn. Syst.*, pp. 1–13, 2021. doi: 10.1109/TNNLS.2021.3069058
- [41] W.-D. Xi, L. Huang, C.-D. Wang, Y. -Y. Zheng and J. -H. Lai, "Deep Rating and Review Neural Network for Item Recommendation," *IEEE Trans. Neural Networks Learn. Syst.*, pp. 1–11, 2021. doi: 10.1109/TNNLS.2021.3083264.
- [42] H. R. Chamorro, I. Riano, R. Gerndt, I. Zelinka, F. Gonzalez-Longatt, and V. K. Sood, "Synthetic inertia control based on fuzzy adaptive differential evolution," *International Journal of Electrical Power & Energy Systems*, vol. 105, pp. 803–813, 2019.
- [43] S. Prabha and R. Yadav, "Differential evolution with biological-based mutation operator," *Engineering Science and Technology, an International Journal*, vol. 23, no. 2, pp. 253–263, 2020.
- [44] S. Li, Q. Gu, W. Gong, and B. Ning, "An enhanced adaptive differential evolution algorithm for parameter extraction of photovoltaic models," *Energy Conversion and Management*, vol. 205, p. 112443, 2020.
- [45] N. B. Guedria, "An accelerated differential evolution algorithm with new operators for multi-damage detection in plate-like structures," *Applied Mathematical Modelling*, vol. 80, pp. 366–383, 2020.
- [46] M. Garanayak, S. N. Mohanty, A. K. Jagadev, and S. Sahoo, "Recommender system using item based collaborative filtering (CF) and K-means," *International Journal of Knowledge-based and Intelligent Engineering Systems*, vol. 23, no.2, pp.93-101, 2019.
- [47] M. Garanayak, G. Sahu, S. N. Mohanty, and A. K. Jagadev, "Agricultural recommendation system for crops using different machine learning regression methods," *International Journal of Agricultural and Environmental Information Systems*, vol. 12, no.1, pp.1-20, 2021.



Di Wu (Member, IEEE) received his Ph.D. degree from the Chongqing Institute of Green and Intelligent Technology (CIGIT), Chinese Academy of Sciences (CAS), China in 2019 and then joined CIGIT, CAS, China. He is currently a Professor of the College of Computer and Information Science, Southwest University, Chongqing, China. He has over 50 publications, including 12 IEEE Transactions papers of IEEE T-KDE, T-NNLS, T-SC, T-SMC, etc., and conferences of AAAI, WWW, ICDM, IJCAI, etc. He is servicing as an Associate Editor for *Frontiers in Neurobotics* (SCI, IF2.65). Homepage: <https://wudi1989.github.io/Homepage/>.



Bo Sun received the B.S. degree in Optoelectronic information science and Engineering from Tianjin University, Tianjin, China, in 2019. He is currently pursuing his M.S. degree in computer technology at University of Chinese Academy of Sciences (UCAS), Beijing, China. His research interests include data mining and deep learning.



Mingsheng Shang received his Ph.D Degree in Computer Science from University of Electronic Science and Technology of China (UESTC). He joined the Chongqing Institute of Green and Intelligent Technology, Chinese Academy of Sciences, Chongqing, China, in 2015, and is currently a Professor of computer science and engineering. He has over 100 publications including IEEE T-KDE, T-NNLS, T-Cyber, T-SMC, T-SC, etc. His research interests include data mining, complex networks, cloud computing, and their applications.